

Using MaxSAT for Efficient Explanations of Tree Ensembles

Alexey Ignatiev¹, Yacine Izza², Peter J. Stuckey¹, Joao Marques-Silva³

February 22 – March 1, 2022 | **AAAI**

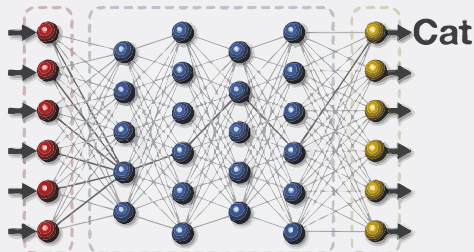
¹Monash University, Melbourne, Australia

²University of Toulouse, France

³IRIT, CNRS, Toulouse, France

Formal eXplainable AI

Machine Learning System



This is a cat.

Current Explanation


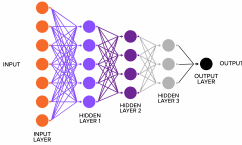






This is a cat:

- It has fur, whiskers, and claws.
- It has this feature:



XAI Explanation

Why? Status quo...

	A parrot	Machine learning algorithm
		
Learns random phrases		
Doesn't understand s**t about what it learns		
Occasionally speaks nonsense		

classifier $\tau : \mathbb{F} \rightarrow \mathcal{K}$, instance \mathbf{v} s.t. $\tau(\mathbf{v}) = \mathbf{c}$

classifier $\tau : \mathbb{F} \rightarrow \mathcal{K}$, instance \mathbf{v} s.t. $\tau(\mathbf{v}) = \mathbf{c}$

abductive explanation \mathcal{X}

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = \mathbf{c})$$

classifier $\tau : \mathbb{F} \rightarrow \mathcal{K}$, instance \mathbf{v} s.t. $\tau(\mathbf{v}) = \mathbf{c}$

abductive explanation \mathcal{X}

“why prediction \mathbf{c} ?”

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = \mathbf{c})$$

classifier $\tau : \mathbb{F} \rightarrow \mathbb{K}$, instance \mathbf{v} s.t. $\tau(\mathbf{v}) = \mathbf{c}$

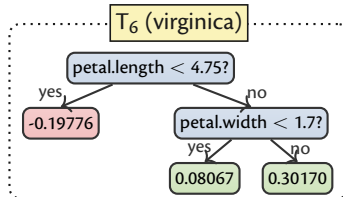
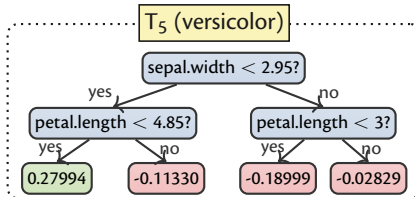
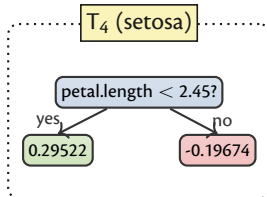
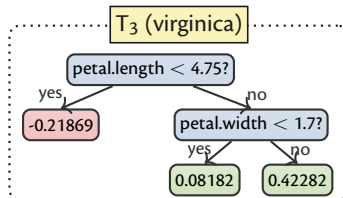
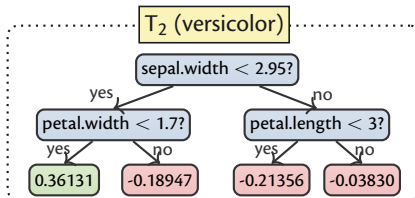
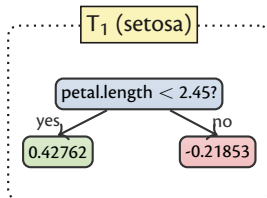
abductive explanation \mathcal{X}

“why prediction \mathbf{c} ?”

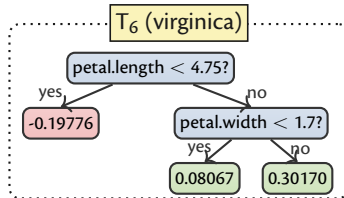
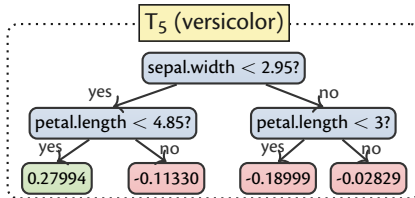
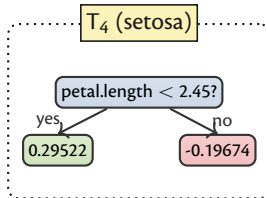
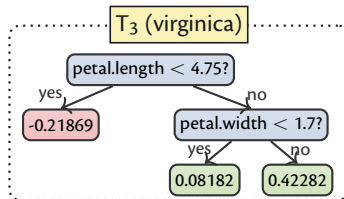
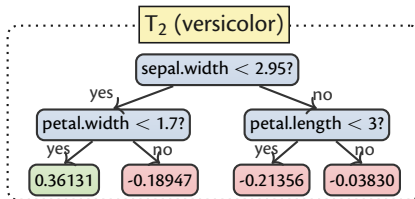
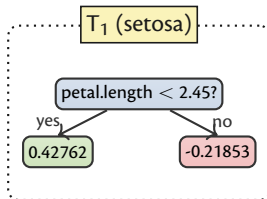
$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = \mathbf{c})$$

because of features of \mathcal{X} !

Abductive explanations for tree ensembles

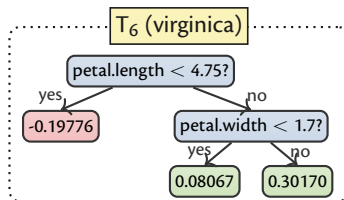
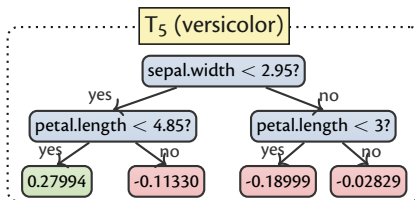
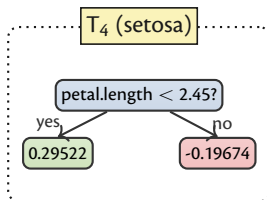
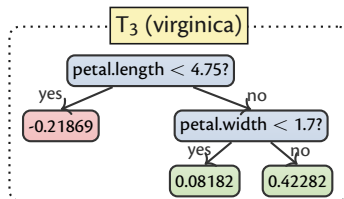
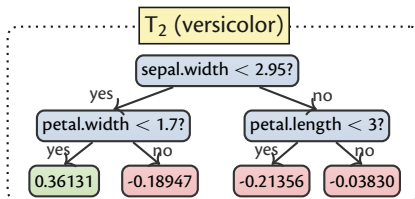
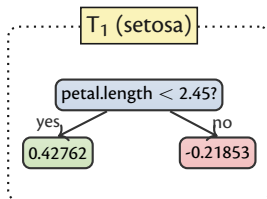


Abductive explanations for tree ensembles



- $w(\mathbf{x}, c) = \sum_{j \in \{0, \dots, n-1\}} \mathcal{T}_{Kj+c}(\mathbf{x}), c \in [K]$
- $\tau(\mathbf{x}) = \arg \max_{c \in [K]} w(\mathbf{x}, c)$

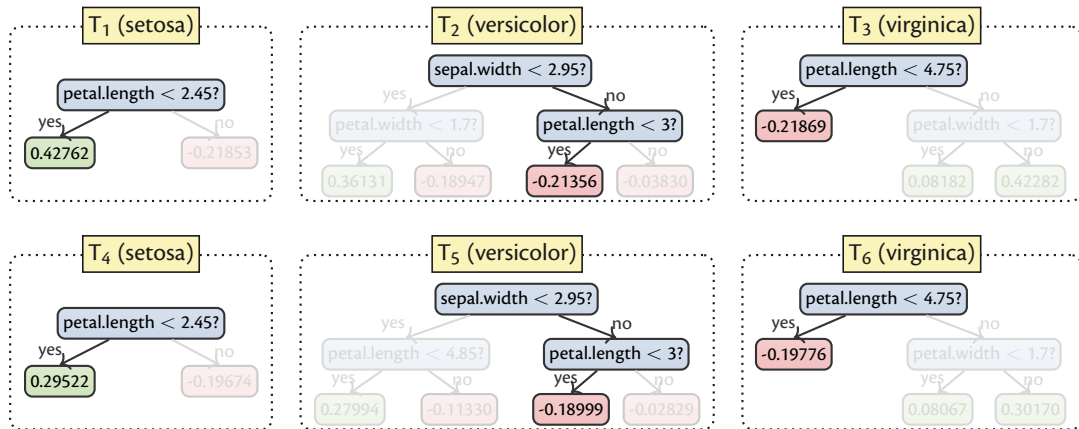
Abductive explanations for tree ensembles



- $w(\mathbf{x}, c) = \sum_{j \in \{0, \dots, n-1\}} \mathcal{T}_{Kj+c}(\mathbf{x}), c \in [K]$
- $\tau(\mathbf{x}) = \arg \max_{c \in [K]} w(\mathbf{x}, c)$

$$\forall (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c)$$

Abductive explanations for tree ensembles



- $w(\mathbf{x}, c) = \sum_{j \in \{0, \dots, n-1\}} \mathcal{T}_{Kj+c}(\mathbf{x}), c \in [K]$
- $\tau(\mathbf{x}) = \arg \max_{c \in [K]} w(\mathbf{x}, c)$

(sepal.length = 5.1) \wedge (sepal.width = 3.5) \wedge (petal.length = 1.4) \wedge (petal.width = 0.2)

$$\forall (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c)$$

Formal explanations with MaxSAT

Overview of the approach

SMT

- reach logic, can handle *linear constraints*

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_i \geq \sum_i \right)$$

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_i \geq \sum_i \right)$$

encode BT operation

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_l \geq \sum_i \right)$$

encode BT operation

enforce class c_l instead of c_i

insane decimal point precision

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_l \geq \sum_i \right)$$

encode BT operation

enforce class c_l instead of c_i

insane decimal point precision



MaxSAT

maximize $\sum_l - \sum_i$
subject to \mathcal{H}

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_l \geq \sum_i \right)$$

encode BT operation

enforce class c_l instead of c_i

insane decimal point precision



MaxSAT

maximize $\sum_l - \sum_i$
subject to \mathcal{H}

- pure propositional logic
 - prediction as objective function
 - weighted soft clauses keep precision
 - core-guided MaxSAT

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_l \geq \sum_i \right)$$

encode BT operation

enforce class c_l instead of c_i

insane decimal point precision



MaxSAT

maximize $\sum_l - \sum_i$
subject to \mathcal{H}

-
- pure propositional logic
 - prediction as objective function
 - weighted soft clauses keep precision
 - core-guided MaxSAT
 - incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_l \geq \sum_i \right)$$

encode BT operation

enforce class c_l instead of c_i

insane decimal point precision



MaxSAT

maximize $\sum_l - \sum_i$
subject to \mathcal{H}

- pure propositional logic
 - prediction as objective function
 - weighted soft clauses keep precision
 - core-guided MaxSAT
- incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse
- novel weight stratification

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_l \geq \sum_i \right)$$

encode BT operation

enforce class c_l instead of c_i

insane decimal point precision



MaxSAT

maximize $\sum_l - \sum_i$
subject to \mathcal{H}

-
- pure propositional logic
 - prediction as objective function
 - weighted soft clauses keep precision
 - core-guided MaxSAT
 - incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse
 - novel weight stratification
 - early termination

Overview of the approach

SMT

- reach logic, can handle *linear constraints*
- directly reason about:

$$\mathcal{H} \wedge \left(\sum_l \geq \sum_i \right)$$

encode BT operation

enforce class c_l instead of c_i

insane decimal point precision



MaxSAT

maximize $\sum_l - \sum_i$
subject to \mathcal{H}

- pure propositional logic
 - prediction as objective function
 - weighted soft clauses keep precision
 - core-guided MaxSAT
- incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse
- novel weight stratification
- early termination

the devil is in the details

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathfrak{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathfrak{E}} \mathcal{H}_{T_i}$$

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathfrak{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathfrak{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathcal{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

- feature threshold values $s_{j,k}$ from \mathcal{E}

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathcal{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

- feature threshold values $s_{j,k}$ from \mathcal{E}
- variable $o_{j,k} = 1$ iff $x_j < s_{j,k}$
- variable $l_{j,k} = 1$ iff $x_j \in [s_{j,k-1}, s_{j,k})$

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathcal{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

- feature threshold values $s_{j,k}$ from \mathcal{E}
- variable $o_{j,k} = 1$ iff $x_j < s_{j,k}$
- variable $l_{j,k} = 1$ iff $x_j \in [s_{j,k-1}, s_{j,k})$
- order encoding of D_j with $l_{j,k}$ and $o_{j,k}$
- instance is expressed by $l_{j,k}$ -variables

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathcal{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

- feature threshold values $s_{j,k}$ from \mathcal{E}
- variable $o_{j,k} = 1$ iff $x_j < s_{j,k}$
- variable $l_{j,k} = 1$ iff $x_j \in [s_{j,k-1}, s_{j,k})$
- order encoding of D_j with $l_{j,k}$ and $o_{j,k}$
- instance is expressed by $l_{j,k}$ -variables

\mathcal{H}_{T_i} encodes paths of T_i

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathcal{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

- feature threshold values $s_{j,k}$ from \mathcal{E}
- variable $o_{j,k} = 1$ iff $x_j < s_{j,k}$
- variable $l_{j,k} = 1$ iff $x_j \in [s_{j,k-1}, s_{j,k})$
- order encoding of D_j with $l_{j,k}$ and $o_{j,k}$
- instance is expressed by $l_{j,k}$ -variables

\mathcal{H}_{T_i} encodes paths of T_i

- \mathcal{P}_i is the set of paths of T_i
- given a $P_r \in \mathcal{P}_i$, t_r denotes its leaf

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathcal{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

- feature threshold values $s_{j,k}$ from \mathcal{E}
- variable $o_{j,k} = 1$ iff $x_j < s_{j,k}$
- variable $l_{j,k} = 1$ iff $x_j \in [s_{j,k-1}, s_{j,k})$
- order encoding of D_j with $l_{j,k}$ and $o_{j,k}$
- instance is expressed by $l_{j,k}$ -variables

\mathcal{H}_{T_i} encodes paths of T_i

- \mathcal{P}_i is the set of paths of T_i
- given a $P_r \in \mathcal{P}_i$, t_r denotes its leaf
- encode each path $P_r \in \mathcal{P}_i$:

$$\left(\bigwedge_{(x_j < s_{j,k}) \in P_r} o_{j,k} \wedge \bigwedge_{(x_j \geq s_{j,k}) \in P_r} \neg o_{j,k} \right) \leftrightarrow t_r$$

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathcal{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

- feature threshold values $s_{j,k}$ from \mathcal{E}
- variable $o_{j,k} = 1$ iff $x_j < s_{j,k}$
- variable $l_{j,k} = 1$ iff $x_j \in [s_{j,k-1}, s_{j,k})$
- order encoding of D_j with $l_{j,k}$ and $o_{j,k}$
- instance is expressed by $l_{j,k}$ -variables

\mathcal{H}_{T_i} encodes paths of T_i

- \mathcal{P}_i is the set of paths of T_i
- given a $P_r \in \mathcal{P}_i$, t_r denotes its leaf
- encode each path $P_r \in \mathcal{P}_i$:
$$\left(\bigwedge_{(x_j < s_{j,k}) \in P_r} o_{j,k} \wedge \bigwedge_{(x_j \geq s_{j,k}) \in P_r} \neg o_{j,k} \right) \leftrightarrow t_r$$
- $\sum_{P_r \in \mathcal{P}_i} t_r = 1$

Encoding BT operation

\mathcal{F} — set of features

$\forall j \in \mathcal{F} \ D_j$ — domain of feature j

\mathcal{E} — tree ensemble

$$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathcal{E}} \mathcal{H}_{T_i}$$

\mathcal{H}_{D_j} encodes feature domain D_j

- feature threshold values $s_{j,k}$ from \mathcal{E}
- variable $o_{j,k} = 1$ iff $x_j < s_{j,k}$
- variable $l_{j,k} = 1$ iff $x_j \in [s_{j,k-1}, s_{j,k})$
- order encoding of D_j with $l_{j,k}$ and $o_{j,k}$
- instance is expressed by $l_{j,k}$ -variables

\mathcal{H}_{T_i} encodes paths of T_i

- \mathcal{P}_i is the set of paths of T_i
- given a $P_r \in \mathcal{P}_i$, t_r denotes its leaf
- encode each path $P_r \in \mathcal{P}_i$:
$$\left(\bigwedge_{(x_j < s_{j,k}) \in P_r} o_{j,k} \wedge \bigwedge_{(x_j \geq s_{j,k}) \in P_r} \neg o_{j,k} \right) \leftrightarrow t_r$$
- $\sum_{P_r \in \mathcal{P}_i} t_r = 1$

see paper for details

“Optimizing” model predictions

AXp condition for \mathcal{X}

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$

“Optimizing” model predictions

AXp condition for \mathcal{X}

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$



\mathcal{X} is not AXp

“Optimizing” model predictions

AXp condition for \mathcal{X}

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$



\mathcal{X} is not AXp if $\exists_{\mathbf{x} \in \mathbb{F}, c_i \neq c_i} \cdot w_l(\mathbf{x}) \geq w_i(\mathbf{x})$

“Optimizing” model predictions

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$

AXp condition for \mathcal{X}



\mathcal{X} is not AXp if $\exists_{\mathbf{x} \in \mathbb{F}, c_l \neq c_i} \cdot w_l(\mathbf{x}) \geq w_i(\mathbf{x})$

weighted sums \sum_l and \sum_i

“Optimizing” model predictions

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$

AXp condition for \mathcal{X}



\mathcal{X} is not AXp if $\exists_{\mathbf{x} \in \mathbb{F}, c_l \neq c_i} \cdot w_l(\mathbf{x}) \geq w_i(\mathbf{x})$

weighted sums \sum_l and \sum_i

hard to reason about

"Optimizing" model predictions

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$

AXp condition for \mathcal{X}



\mathcal{X} is not AXp if $\exists_{\mathbf{x} \in \mathbb{F}, c_i \neq c_i} \cdot w_l(\mathbf{x}) \geq w_i(\mathbf{x})$

weighted sums \sum_l and \sum_i

even to represent in propositional logic!

hard to reason about

“Optimizing” model predictions

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$

AXp condition for \mathcal{X}



\mathcal{X} is not AXp if $\exists_{\mathbf{x} \in \mathbb{F}, c_i \neq c_i} \cdot w_l(\mathbf{x}) \geq w_i(\mathbf{x})$

weighted sums \sum_l and \sum_i

even to represent in propositional logic!

hard to reason about



maximize

subject to

$$s_{i,l} = \sum_l - \sum_i$$
$$\mathcal{H} \wedge \bigwedge_{j \in \mathcal{X}} \mathbb{I}[x_j = v_j]$$

"Optimizing" model predictions

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$

AXp condition for \mathcal{X}



\mathcal{X} is not AXp if $\exists_{\mathbf{x} \in \mathbb{F}, c_i \neq c_i} \cdot w_l(\mathbf{x}) \geq w_i(\mathbf{x})$

weighted sums \sum_l and \sum_i

even to represent in propositional logic!

hard to reason about



\mathcal{X} is not AXp if $S_{i,l}^* \geq 0$

maximize

subject to

$$S_{i,l} = \sum_l - \sum_i$$

$$\mathcal{H} \wedge \bigwedge_{j \in \mathcal{X}} \mathbb{I}[x_j = v_j]$$

“Optimizing” model predictions

$$\forall (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c_i)$$

AXp condition for \mathcal{X}



\mathcal{X} is not AXp if $\exists_{\mathbf{x} \in \mathbb{F}, c_i \neq c_i} \cdot w_l(\mathbf{x}) \geq w_i(\mathbf{x})$

weighted sums \sum_l and \sum_i

even to represent in propositional logic!

hard to reason about



\mathcal{X} is not AXp if $S_{i,l}^* \geq 0$

maximize

subject to

$$S_{i,l} = \sum_l - \sum_i$$

$$\mathcal{H} \wedge \bigwedge_{j \in \mathcal{X}} \mathbb{I}[x_j = v_j]$$

weighted soft clauses + propositional variables only

multiple calls to a MaxSAT solver
with varying assumptions \mathcal{X}

multiple calls to a MaxSAT solver
with varying assumptions \mathcal{X}



the need for incremental MaxSAT!

Incrementality in MaxSAT

Function MaxSAT(ϕ, \mathcal{A}):

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost  $\leftarrow$  0                                # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$     # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                             # iterate over known cores  $\kappa$ :
4      cost  $\leftarrow$  cost + COREWT( $\kappa$ ) # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$  # process  $\kappa$  and update  $\phi$ 
6  while SAT( $\phi, \mathcal{A}$ ) = false: # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$            # new unsatisfiable core
8      cost  $\leftarrow$  cost + COREWT( $\kappa$ ) # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$  # process  $\kappa$  and update  $\phi$ 
10     RECORD( $\phi, \mathcal{A}, \kappa$ )                # record  $\kappa$  for the future
11  return GETMODEL( $\phi$ )                        #  $\phi$  is now satisfiable
```


Incrementality in MaxSAT

Function $\text{MaxSAT}(\phi, \mathcal{A})$:

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost ← 0                                # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$     # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                          # iterate over known cores  $\kappa$ 
4      cost ← cost +  $\text{COREWT}(\kappa)$  # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$     # process  $\kappa$  and update  $\phi$ 
6  while  $\text{SAT}(\phi, \mathcal{A}) = \text{false}$ : # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$           # new unsatisfiable core
8      cost ← cost +  $\text{COREWT}(\kappa)$  # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$     # process  $\kappa$  and update  $\phi$ 
10      $\text{RECORD}(\phi, \mathcal{A}, \kappa)$           # record  $\kappa$  for the future
11  return  $\text{GETMODEL}(\phi)$                 #  $\phi$  is now satisfiable
```

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{\perp[x_j = v_j] \mid j \in \mathcal{X}\}$

Incrementality in MaxSAT

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{\llbracket x_j = v_j \rrbracket \mid j \in \mathcal{X}\}$

Function MaxSAT(ϕ, \mathcal{A}):

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost ← 0                                # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$   # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                        # iterate over known cores  $\kappa$ 
4      assumptions are passed down to SAT solver
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$       # process  $\kappa$  and update  $\phi$ 
6  while SAT( $\phi, \mathcal{A}$ ) = false: # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$           # new unsatisfiable core
8      cost ← cost + COREWT( $\kappa$ ) # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$       # process  $\kappa$  and update  $\phi$ 
10     RECORD( $\phi, \mathcal{A}, \kappa$ )              # record  $\kappa$  for the future
11  return GETMODEL( $\phi$ )                  #  $\phi$  is now satisfiable
```

Incrementality in MaxSAT

Function $\text{MaxSAT}(\phi, \mathcal{A})$:

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost ← 0                                # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$   # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                        # iterate over known cores  $\kappa$ 
4       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$       # process  $\kappa$  and update  $\phi$ 
5       $\text{cost} \leftarrow \text{cost} + \text{COREWT}(\kappa)$  # add its weight to cost
6  while  $\text{SAT}(\phi, \mathcal{A}) = \text{false}$ : # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$           # new unsatisfiable core
8       $\text{cost} \leftarrow \text{cost} + \text{COREWT}(\kappa)$  # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$       # process  $\kappa$  and update  $\phi$ 
10      $\text{RECORD}(\phi, \mathcal{A}, \kappa)$            # record  $\kappa$  for the future
11  return  $\text{GETMODEL}(\phi)$                  #  $\phi$  is now satisfiable
```

assumptions are passed down to SAT solver

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{\perp[x_j = v_j] \mid j \in \mathcal{X}\}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$

Incrementality in MaxSAT

Function $\text{MaxSAT}(\phi, \mathcal{A})$:

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost ← 0                                # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$   # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                        # iterate over known cores  $\kappa$ 
4       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$     # process  $\kappa$  and update  $\phi$ 
5       $\text{cost} \leftarrow \text{cost} + \text{COREWT}(\kappa)$  # add its weight to cost
6  while  $\text{SAT}(\phi, \mathcal{A}) = \text{false}$ : # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$          # new unsatisfiable core
8       $\text{cost} \leftarrow \text{cost} + \text{COREWT}(\kappa)$  # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$     # process  $\kappa$  and update  $\phi$ 
10      $\text{RECORD}(\phi, \mathcal{A}, \kappa)$           # record  $\kappa$  for the future
11  return  $\text{GETMODEL}(\phi)$                 #  $\phi$  is now satisfiable
```

assumptions are passed down to SAT solver

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{\llbracket x_j = v_j \rrbracket \mid j \in \mathcal{X}\}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- *similar to modern SAT solvers*

Incrementality in MaxSAT

Function $\text{MaxSAT}(\phi, \mathcal{A})$:

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost ← 0                                # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$     # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                          # iterate over known cores  $\kappa$ 
4      cost ← cost +  $\text{COREWT}(\kappa)$  # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$     # process  $\kappa$  and update  $\phi$ 
6  while  $\text{SAT}(\phi, \mathcal{A}) = \text{false}$ : # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$           # new unsatisfiable core
8      cost ← cost +  $\text{COREWT}(\kappa)$  # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$     # process  $\kappa$  and update  $\phi$ 
10      $\text{RECORD}(\phi, \mathcal{A}, \kappa)$           # record  $\kappa$  for the future
11 return  $\text{GETMODEL}(\phi)$                   #  $\phi$  is now satisfiable
```

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{\llbracket x_j = v_j \rrbracket \mid j \in \mathcal{X}\}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- *similar to modern SAT solvers*

unsatisfiable core reuse:

Incrementality in MaxSAT

Function MaxSAT(ϕ, \mathcal{A}):

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost  $\leftarrow$  0                                # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$     # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                          # iterate over known cores  $\kappa$ 
4      cost  $\leftarrow$  cost + COREWT( $\kappa$ ) # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$         # process  $\kappa$  and update  $\phi$ 
6  while SAT( $\phi, \mathcal{A}$ ) = false: # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$             # new unsatisfiable core
8      cost  $\leftarrow$  cost + COREWT( $\kappa$ ) # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$         # process  $\kappa$  and update  $\phi$ 
10 RECORD( $\phi, \mathcal{A}, \kappa$ )                      # record  $\kappa$  for the future
11 return GETMODEL( $\phi$ )                        #  $\phi$  is now satisfiable
```

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{\perp[x_j = v_j] \mid j \in \mathcal{X}\}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- *similar to modern SAT solvers*

unsatisfiable core reuse:

- record each new core κ_i

Incrementality in MaxSAT

Function MaxSAT(ϕ, \mathcal{A}):

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost  $\leftarrow$  0                                # initially, cost is 0

2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$       # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                             # iterate over known cores  $\kappa$ 
4      cost  $\leftarrow$  cost + COREWT( $\kappa$ ) # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$          # process  $\kappa$  and update  $\phi$ 

6  while SAT( $\phi, \mathcal{A}$ ) = false: # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$               # new unsatisfiable core
8      cost  $\leftarrow$  cost + COREWT( $\kappa$ ) # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$          # process  $\kappa$  and update  $\phi$ 
10     RECORD( $\phi, \mathcal{A}, \kappa$ )                  # record  $\kappa$  for the future
11 return GETMODEL( $\phi$ )                          #  $\phi$  is now satisfiable
```

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{ \perp[x_j = v_j] \mid j \in \mathcal{X} \}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- *similar to modern SAT solvers*

unsatisfiable core reuse:

- record each new core κ_i
- start a new call by reusing known cores

Incrementality in MaxSAT

Function MaxSAT(ϕ, \mathcal{A}):

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost  $\leftarrow$  0                                # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$     # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                          # iterate over known cores  $\kappa$ 
4      cost  $\leftarrow$  cost + COREWT( $\kappa$ ) # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$         # process  $\kappa$  and update  $\phi$ 
6  while SAT( $\phi, \mathcal{A}$ ) = false: # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$             # new unsatisfiable core
8      cost  $\leftarrow$  cost + COREWT( $\kappa$ ) # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$         # process  $\kappa$  and update  $\phi$ 
10     RECORD( $\phi, \mathcal{A}, \kappa$ )                # record  $\kappa$  for the future
11 return GETMODEL( $\phi$ )                        #  $\phi$  is now satisfiable
```

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{\perp[x_j = v_j] \mid j \in \mathcal{X}\}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- similar to modern SAT solvers*

unsatisfiable core reuse:

- record each new core κ_i
 - use external SAT solver
 - record all the literals core κ_i depends on
 - i.e. use $\bigwedge_{l \in \mathcal{D}_i} \rightarrow \zeta_i$
- start a new call by **reusing known cores**

Incrementality in MaxSAT

Function $\text{MaxSAT}(\phi, \mathcal{A})$:

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost  $\leftarrow$  0                                     # initially, cost is 0

2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$            # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                                # iterate over known cores  $\kappa$ 
4      cost  $\leftarrow$  cost +  $\text{COREWT}(\kappa)$  # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$            # process  $\kappa$  and update  $\phi$ 

6  while  $\text{SAT}(\phi, \mathcal{A}) = \text{false}$ : # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$                  # new unsatisfiable core
8      cost  $\leftarrow$  cost +  $\text{COREWT}(\kappa)$  # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$            # process  $\kappa$  and update  $\phi$ 
10      $\text{RECORD}(\phi, \mathcal{A}, \kappa)$                  # record  $\kappa$  for the future
11 return  $\text{GETMODEL}(\phi)$                          #  $\phi$  is now satisfiable
```

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{\perp[x_j = v_j] \mid j \in \mathcal{X}\}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- *similar to modern SAT solvers*

unsatisfiable core reuse:

- record each new core κ_i
 - use external SAT solver
 - record all the literals core κ_i depends on
 - i.e. use $\bigwedge_{l \in \mathcal{D}_i} \rightarrow \zeta_i$
- start a new call by **reusing known cores**
 - apply **unit propagation** given \mathcal{S} and \mathcal{A}
 - all reusable cores are “*propagated*” in the right order

Incrementality in MaxSAT

Function MaxSAT(ϕ, \mathcal{A}):

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost ← 0                                     # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$     # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                             # iterate over known cores  $\kappa$ 
4      cost ← cost + COREWT( $\kappa$ ) # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$          # process  $\kappa$  and update  $\phi$ 
6  while SAT( $\phi, \mathcal{A}$ ) = false: # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$              # new unsatisfiable core
8      cost ← cost + COREWT( $\kappa$ ) # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$          # process  $\kappa$  and update  $\phi$ 
10     RECORD( $\phi, \mathcal{A}, \kappa$ )                 # record  $\kappa$  for the future
11  return GETMODEL( $\phi$ )                         #  $\phi$  is now satisfiable
```

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{l[x_j = v_j] \mid j \in \mathcal{X}\}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- similar to modern SAT solvers*

unsatisfiable core reuse:

- record each new core κ_i
 - use external SAT solver
 - record all the literals core κ_i depends on
 - i.e. use $\bigwedge_{l \in \mathcal{D}_i} \rightarrow \zeta_i$
- start a new call by **reusing known cores**
 - apply **unit propagation** given \mathcal{S} and \mathcal{A}
 - all reusable cores are “*propagated*” in the right order

early termination:

- no need to solve the problem to optimality!

Incrementality in MaxSAT

Function MaxSAT(ϕ, \mathcal{A}):

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost ← 0                                     # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$     # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                             # iterate over known cores  $\kappa$ 
4      cost ← cost + COREWT( $\kappa$ ) # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$  # process  $\kappa$  and update  $\phi$ 
6  while SAT( $\phi, \mathcal{A}$ ) = false: # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$  # new unsatisfiable core
8      cost ← cost + COREWT( $\kappa$ ) # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$ 
10
11 return GETMODEL( $\phi$ ) #  $\phi$  is now satisfiable
```

LB on the cost is updated at each iteration

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{ \llbracket x_j = v_j \rrbracket \mid j \in \mathcal{X} \}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- similar to modern SAT solvers*

unsatisfiable core reuse:

- record each new core κ_i
 - use external SAT solver
 - record all the literals core κ_i depends on
 - i.e. use $\bigwedge_{l \in \mathcal{D}_i} \rightarrow \zeta_i$
- start a new call by **reusing known cores**
 - apply **unit propagation** given \mathcal{S} and \mathcal{A}
 - all reusable cores are “*propagated*” in the right order

early termination:

- no need to solve the problem to optimality!
- no misclassification** if $\text{cost} > \sum_i$

Incrementality in MaxSAT

Function MaxSAT(ϕ, \mathcal{A}):

Input: ϕ : Partial CNF formula ($\phi \triangleq \mathcal{H} \wedge \mathcal{S}$)

\mathcal{A} : Set of assumption literals

Output: μ : MaxSAT model

```
1  cost ← 0                                     # initially, cost is 0
2   $\mathcal{C} \leftarrow \text{VALIDCORES}(\phi, \mathcal{A})$     # get valid unsatisfiable cores
3  foreach  $\kappa \in \mathcal{C}$ :                             # iterate over known cores  $\kappa$ 
4      cost ← cost + COREWT( $\kappa$ ) # add its weight to cost
5       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$  # process  $\kappa$  and update  $\phi$ 
6  while SAT( $\phi, \mathcal{A}$ ) = false: # iterate until  $\phi$  gets satisfiable
7       $\kappa \leftarrow \text{GETCORE}(\phi)$  # new unsatisfiable core
8      cost ← cost + COREWT( $\kappa$ ) # add its weight to cost
9       $\phi \leftarrow \text{PROCESS}(\phi, \kappa)$  # process  $\kappa$  and update  $\phi$ 
10     RECORD( $\phi, \mathcal{A}, \kappa$ ) # record  $\kappa$  for the future
11  return GETMODEL( $\phi$ ) #  $\phi$  is now satisfiable
```

MiniSat-like assumptions interface:

- assumptions $\mathcal{A} \triangleq \{ \perp[x_j = v_j] \mid j \in \mathcal{X} \}$
- can get “*unsat core*” $\mathcal{A}^* = \cup \mathcal{A}'$ s.t. $\mathcal{A}' = \kappa \cap \mathcal{A}$
- similar to modern SAT solvers*

unsatisfiable core reuse:

- record each new core κ_i
 - use external SAT solver
 - record all the literals core κ_i depends on
 - i.e. use $\bigwedge_{l \in \mathcal{D}_i} \rightarrow \zeta_i$
- start a new call by **reusing known cores**
 - apply **unit propagation** given \mathcal{S} and \mathcal{A}
 - all reusable cores are “*propagated*” in the right order

early termination:

- no need to solve the problem to optimality!
- no misclassification** if $\text{cost} > \sum_i$
- misclassification** if over-approximation's objective > 0

can over-approximate objective once current level is solved

Additional heuristic — distance-based stratification

why?

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$

$(t_2, 0.41527)$

$(\neg t_3, 0.41645)$

$(t_4, 0.16249)$

$(t_5, 0.72452)$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$\mathcal{L}_1 = \left\{ (t_5, 0.72452) \right\}$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$0.00168 = |0.72284 - 0.72452| < |0.72284 - (0.41645 + 0.41527 + 0.16249)/3| = 0.39144$$

$$\mathcal{L}_1 = \left\{ (t_5, 0.72452) \right\}$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$\mathcal{L}_1 = \left\{ \begin{array}{l} (t_5, 0.72452) \\ (\neg t_1, 0.72284) \end{array} \right\}$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$0.30723 = |0.41645 - 0.72452 + 0.72284/2| > |0.41645 - (0.41527 + 0.16249)/2| = 0.12757$$

$$\mathcal{L}_1 = \left\{ \begin{array}{l} (t_5, 0.72452) \\ (\neg t_1, 0.72284) \end{array} \right\}$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$\mathcal{L}_1 = \left\{ \begin{array}{l} (t_5, 0.72452) \\ (\neg t_1, 0.72284) \end{array} \right\}$$

$$\mathcal{L}_2 = \left\{ (\neg t_3, 0.41645) \right\}$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$0.00118 = |0.41527 - 0.41645| < |0.41527 - 0.16249| = 0.25278$$

$$\mathcal{L}_1 = \left\{ \begin{array}{l} (t_5, 0.72452) \\ (\neg t_1, 0.72284) \end{array} \right\}$$

$$\mathcal{L}_2 = \left\{ (\neg t_3, 0.41645) \right\}$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$\mathcal{L}_1 = \left\{ \begin{array}{l} (t_5, 0.72452) \\ (\neg t_1, 0.72284) \end{array} \right\}$$

$$\mathcal{L}_2 = \left\{ \begin{array}{l} (\neg t_3, 0.41645) \\ (t_2, 0.41527) \end{array} \right\}$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$0.25337 = |0.16249 - 0.41527 + 0.41645/2| > |0.16249 - 0| = 0.16249$$

$$\mathcal{L}_1 = \left\{ \begin{array}{l} (t_5, 0.72452) \\ (\neg t_1, 0.72284) \end{array} \right\} \quad \mathcal{L}_2 = \left\{ \begin{array}{l} (\neg t_3, 0.41645) \\ (t_2, 0.41527) \end{array} \right\}$$

Additional heuristic — distance-based stratification

why? — many weighted clauses with **unique weights**

BLO and diversity-based stratification fail to apply!



$$\left| w - \frac{\sum_{(c_i, w_i) \in \mathcal{L}} w_i}{|\mathcal{L}|} \right| < \left| w - \frac{\sum_{(c_j, w_j) \in \mathcal{S} \wedge w_j < w} w_j}{|\{(c_j, w_j) \in \mathcal{S} \wedge w_j < w\}|} \right|$$

$(\neg t_1, 0.72284)$ $(t_2, 0.41527)$ $(\neg t_3, 0.41645)$ $(t_4, 0.16249)$ $(t_5, 0.72452)$

$$\mathcal{L}_1 = \left\{ \begin{array}{l} (t_5, 0.72452) \\ (\neg t_1, 0.72284) \end{array} \right\}$$

$$\mathcal{L}_2 = \left\{ \begin{array}{l} (\neg t_3, 0.41645) \\ (t_2, 0.41527) \end{array} \right\}$$

$$\mathcal{L}_3 = \{(t_4, 0.16249)\}$$

Experimental results

Experimental setup

- **machine configuration:**
 - Intel Core i5 2.30GHz with 8GByte RAM

Experimental setup

- **machine configuration:**

- Intel Core i5 2.30GHz with 8GByte RAM
- running macOS Big Sur

Experimental setup

- **machine configuration:**

- Intel Core i5 2.30GHz with 8GByte RAM
- running macOS Big Sur
- no time or memory limit

Experimental setup

- **machine configuration:**

- Intel Core i5 2.30GHz with 8GByte RAM
- running macOS Big Sur
- no time or memory limit

- **UCI Machine Learning Repository** + **Penn Machine Learning Benchmarks**

Experimental setup

- **machine configuration:**

- Intel Core i5 2.30GHz with 8GByte RAM
- running macOS Big Sur
- no time or memory limit

- **UCI Machine Learning Repository** + **Penn Machine Learning Benchmarks**

- **21 datasets**

- 7–60 features, 26.95 on average
- 59–58000 data instances, 4684.47 on average

Experimental setup

- **machine configuration:**

- Intel Core i5 2.30GHz with 8GByte RAM
- running macOS Big Sur
- no time or memory limit

- **UCI Machine Learning Repository** + **Penn Machine Learning Benchmarks**

- **21 datasets**

- 7–60 features, 26.95 on average
- 59–58000 data instances, 4684.47 on average

- **XGBoost models**

- depth 3–4
- 50 trees per class

Experimental setup

- **machine configuration:**

- Intel Core i5 2.30GHz with 8GByte RAM
- running macOS Big Sur
- no time or memory limit

- **UCI Machine Learning Repository** + **Penn Machine Learning Benchmarks**

- **21 datasets**

- 7–60 features, 26.95 on average
- 59–58000 data instances, 4684.47 on average

- **XGBoost models**

- depth 3–4
- 50 trees per class

- **200 randomly selected** instances to explain

- 3665 individual benchmarks in total

Experimental setup

- **competition tested:**
 - **SMT-based formal explainer**
 - makes calls to **Z3** through **PySMT**
 - same explanation quality

¹<https://github.com/alexeyignatiev/xreason/>

- **competition tested:**
 - **SMT-based formal explainer**
 - makes calls to **Z3** through **PySMT**
 - same explanation quality
 - **Anchor heuristic explainer**
 - sampling-based
 - *no* explanation quality guarantees
 - all parameters are set to default values

¹<https://github.com/alexeyignatiev/xreason/>

Experimental setup

- **competition tested:**

- **SMT-based formal explainer**

- makes calls to **Z3** through **PySMT**
 - same explanation quality

- **Anchor heuristic explainer**

- sampling-based
 - *no* explanation quality guarantees
 - all parameters are set to default values

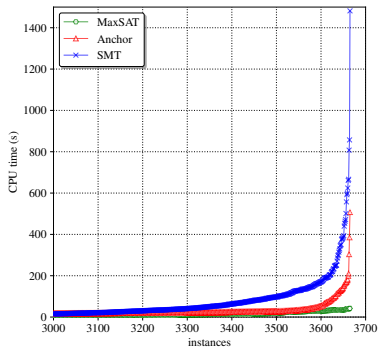
- **prototype¹**

- **MaxSAT-based**

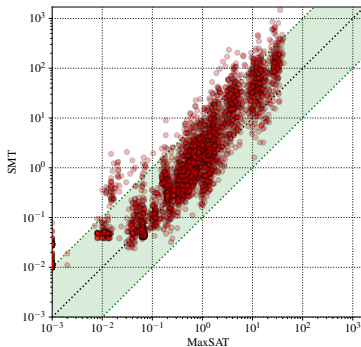
- **same code base** as in the SMT-based variant
 - applies **incrementally extended RC2**
 - makes calls to **Glucose 3** through **PySAT**

¹<https://github.com/alexeyignatiev/xreason/>

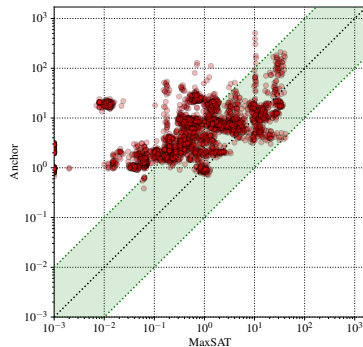
Experimental results



all three competitors

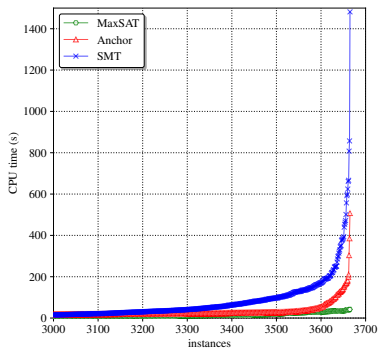


MaxSAT vs. SMT

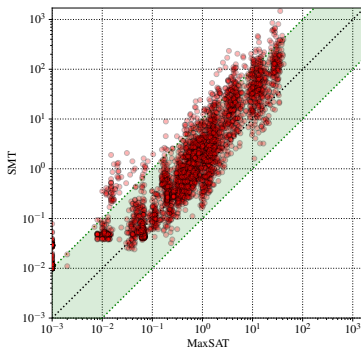


MaxSAT vs. Anchor

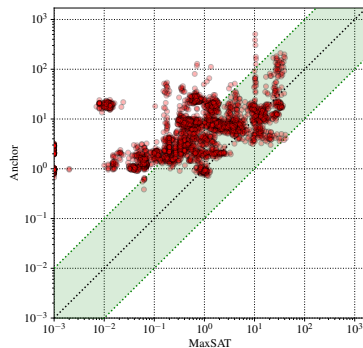
Experimental results



all three competitors



MaxSAT vs. SMT



MaxSAT vs. Anchor

up to 2 orders of magnitude performance improvement

more robust than SMT and Anchor

(see paper for details)

Summary and future work

Summary and future work

- **novel approach** to explaining boosted trees
 - pure propositional logic
 - prediction as objective function
 - core-guided MaxSAT

Summary and future work

- **novel approach** to explaining boosted trees
 - pure propositional logic
 - prediction as objective function
 - core-guided MaxSAT
 - incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse

Summary and future work

- **novel approach** to explaining boosted trees
 - pure propositional logic
 - prediction as objective function
 - core-guided MaxSAT
 - incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse
 - novel weight stratification
 - early termination

Summary and future work

- **novel approach** to explaining boosted trees
 - pure propositional logic
 - prediction as objective function
 - core-guided MaxSAT
 - incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse
 - novel weight stratification
 - early termination
 - a few orders of magnitude performance improvement

Summary and future work

- **novel approach** to explaining boosted trees

- pure propositional logic
 - prediction as objective function
 - core-guided MaxSAT
- incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse
- novel weight stratification
- early termination
- a few orders of magnitude performance improvement

- future work

- contrastive explanations with MaxSAT

Summary and future work

- **novel approach** to explaining boosted trees
 - pure propositional logic
 - prediction as objective function
 - core-guided MaxSAT
 - incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse
 - novel weight stratification
 - early termination
 - a few orders of magnitude performance improvement
- future work
 - contrastive explanations with MaxSAT
 - same ideas for other ML models?

Summary and future work

- **novel approach** to explaining boosted trees
 - pure propositional logic
 - prediction as objective function
 - core-guided MaxSAT
 - incremental MaxSAT calls
 - *MiniSat-like assumptions interface!*
 - unsatisfiable core reuse
 - novel weight stratification
 - early termination
 - a few orders of magnitude performance improvement
- future work
 - contrastive explanations with MaxSAT
 - same ideas for other ML models?
 - more applications of incremental MaxSAT?

Questions?